


RESEARCH ARTICLE

Improving phase II oncology trials using best observed RECIST response as an endpoint by modelling continuous tumour measurements

Chien-Ju Lin¹  | James M.S. Wason^{1,2}¹MRC Biostatistics Unit, University of Cambridge, U.K.²Institute of Health and Society, Newcastle University, U.K.**Correspondence**

Chien-Ju Lin, MRC Biostatistics Unit,
Cambridge Institute of Public Health,
Robinson Way, Cambridge CB2 0SR, U.K.
Email: chienju@mrc-bsu.cam.ac.uk

Funding information

Medical Research Council, Grant/Award
Number: MC_UP_1302/4; Cancer
Research UK, Grant/Award Number:
C48553/A18113

In many phase II trials in solid tumours, patients are assessed using endpoints based on the Response Evaluation Criteria in Solid Tumours (RECIST) scale. Often, analyses are based on the response rate. This is the proportion of patients who have an observed tumour shrinkage above a predefined level and no new tumour lesions. The augmented binary method has been proposed to improve the precision of the estimator of the response rate. The method involves modelling the tumour shrinkage to avoid dichotomising it. However, in many trials the best observed response is used as the primary outcome. In such trials, patients are followed until progression, and their best observed RECIST outcome is used as the primary endpoint. In this paper, we propose a method that extends the augmented binary method so that it can be used when the outcome is best observed response. We show through simulated data and data from a real phase II cancer trial that this method improves power in both single-arm and randomised trials. The average gain in power compared to the traditional analysis is equivalent to approximately a 35% increase in sample size. A modified version of the method is proposed to reduce the computational effort required. We show this modified method maintains much of the efficiency advantages.

KEYWORDS

continuous tumour shrinkage endpoints, longitudinal model, phase II cancer trial

1 | INTRODUCTION

A new cancer treatment is tested for potential benefit in phase II trials that use a relatively small number of patients followed over a short period of time. The results of the phase II trial determines whether to test the treatment in a larger, more time-consuming, and more costly phase III trial. Because of the high cost of, and high failure rate in, phase III oncology trials,¹ it is important to improve the analysis of phase II trials to ensure the decision is more accurate.

Phase II oncology trials use a variety of endpoints to evaluate the efficacy of a treatment.^{2,3} The most commonly used endpoints in solid tumours are based on the Response Evaluation Criteria in Solid Tumours (RECIST) scale.⁴ The RECIST defines tumour size as the sum of longest diameters of target lesions and categorises patients into complete response

The copyright line for this article was changed on 20 November 2017 after original online publication.

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2017 The Authors. *Statistics in Medicine* Published by John Wiley & Sons Ltd.

(CR), partial response (PR), stable disease (SD), and progressive disease (PD). The CR and PR represent no new tumour lesions and a 100% shrinkage and greater than 30% shrinkage, respectively; PD represents a 20% increase in tumour size from the minimum size observed up to that point or new lesions appearing. Often patients are followed until they are categorised as PD or a preplanned time, and patients with CR or PR are labelled responders. The endpoints for evaluating a treatment include response rate and time to an event (progression/death). The response rate is defined as either (1) proportion of patients who are responders at a certain time after baseline (fixed response) or (2) the proportion of patients whose best observed response (BOR) before progression is CR or PR. A response endpoint classing CR and PR as success is called objective response rate (ORR), which is often used in practice. Disease control rate additionally includes SD as success. In some cases where response rate may not be the optimal method, a time to event endpoint “progression-free survival” is considered, which is the time until PD or death. Response-based endpoints are an indicator of the relative anti-tumour activity of the treatments and are not always highly predictive of overall survival. Improving the efficiency of trials using response-based endpoints would help increase the success rate of phase III trials in tumour types where response is predictive of overall survival.

Categorizing patients into responders and nonresponders is widespread and clinically appealing. However, it can have substantial statistical disadvantages. Its major limitation is that it dichotomises the continuous tumour variable, thus discarding information. This loses substantial efficiency.⁵ Some researchers have addressed the problem and proposed methods to make use of the continuous change in tumour response to improve statistical efficiency. This is done in different ways. Karrison et al⁶ propose directly using the change in tumour size as an endpoint. Wason and Seaman⁷ use models for the tumour size and new lesion data which can be used to infer the fixed response rate with higher precision. They find the Karrison's method was more powerful when the probability of new lesions was different between arms but less powerful when the mean tumour size change was different. Jaki et al⁸ propose a method that links tumour size change with mortality using historical datasets. Authors have demonstrated that using continuous scales can increase the power (or reduce the required sample for a target power) compared to analysing the binary composite outcome.

Focusing on response-based endpoints, the method of Wason and Seaman⁷ retains the clinically meaningful endpoint but takes into account the continuous information on tumour size. This method is limited by only allowing 2 follow-up visits and only considering response rate at a fixed time (ie, it cannot be used to make inferences on BOR). In trials in which patients are assessed twice (interim and final), their method is sufficient. However, in trials where patients are followed up until a preplanned time, a method incorporating information on all measurement data is preferred for efficiency reasons. In this paper, we consider ORR and present an extended method that can be used for any number of follow-up times for fixed response or BOR. We propose a modified version that uses a highly efficient technique for multivariate integration,⁹ which substantially reduces the computation time taken. We assess the properties of the proposed methods by using simulated data and data from a real phase II cancer trial (HORIZON II).

This paper is divided into 4 sections. Section 2 gives a brief overview of the augmented binary method (Augbin).⁷ It then describes the proposed extensions of the method. Section 3 evaluates the performance of the proposed methods using simulations and real data. Section 4 summarizes the results and presents limitations and future work.

2 | METHODS

2.1 | Background

We use the phrase “tumour size” as shorthand for the sum of the longest diameter of target tumour lesions. We assume patients tumour sizes are recorded until progression occurs or until a preplanned number of visits. We note there are 2 ways in which a progression can occur: an increase in tumour size by more than 20% (a tumour-growth progression) or new lesions appearing (a new-lesion progression). Two response-based endpoints can be used in the analysis, one being *fixed time* and the other being BOR. Analysis at a fixed time t uses the proportion of responders at time t (those who have a tumour size shrinkage at time t above a predefined threshold and no progression up to that point). Best observed response defines patients as a responder or not according to their BOR before progression. The latest RECIST guidelines⁴ give BOR 2 definitions according to whether confirmation is required or not. Confirmation means that an apparent response must be backed up by continued response at the next time point to be counted as genuine. This is especially recommended for single-arm trials. When confirmation is not required (randomised trials comparing 2 arms), BOR is defined as the best response across all time points up to progression. When confirmation is required, BOR is defined as a response if the patient is a responder at 2 consecutive time points before progression.

2.2 | Notation

Tumour sizes for each patient are measured at several discrete times (T denoting the maximum time). The tumour size at time t for patient i is denoted by z_{it} where $t = 0$ represents the baseline measurement. We denote G and X as the time at which a tumour-growth progression and new-lesion progression occurs, respectively. Once a patient progresses they are no longer followed up. The observed data is therefore (F_i, \mathbf{z}_i) where $F_i = \min(X_i, G_i, T)$. We define y_{it} as the log tumour size ratio for patient i at time t , $y_{it} = \log(z_{it}/z_{i0})$, and c_t as the prespecified dichotomisation threshold for response (on the log tumour ratio scale). Further, D_{it} defines new-lesion progression indicators: $\{D_{it} = 1 \text{ if patient } i \text{ has a progression due to new lesions occurring between time } (t-1) \text{ and } t, t = 1, \dots, T\}$. For simplicity, we define composite indicators S for fixed time and BOR for best observed response using tumour progression relative to the baseline rather than nadir (the lowest tumour size observed so far). Note that we show how to use the true definition in the Supporting Information. In addition, we assume that response at a fixed time would mean just the response status at the specified time is of interest. The response indicator for patient i using fixed time is defined as

$$S_{it} = \begin{cases} 1, & \text{if } D_{ij} = 0 \text{ for all } j = 1, \dots, t, y_{it} < c_t \text{ and no tumour progression before } t, \\ 0, & \text{otherwise.} \end{cases}$$

For BOR, when confirmation is not required, the event is equivalent to having at least 1 record classified CR/PR before progression or time T , the response indicator BOR_i is defined as

$$BOR_i = \begin{cases} 1, & \text{if there exists a } t \text{ such that } y_{it} < \log(0.7), t \leq \min(F_i, X_i - 1) \text{ and } \log(0.7) < y_{i1}, \dots, y_{i(t-1)} < \log(1.2), \\ 0, & \text{otherwise.} \end{cases}$$

We consider the case where confirmation is required later. We lay out more fully how to incorporate the actual definition of progression with change from nadir in the Supporting Information. This includes use of an indicator function for the extended augmented binary method (eAugbin) and a more exact method for the modified augmented binary method (mAug).

2.3 | Estimating response probability using the Augbin with 2 follow-up times

The augmented binary method, henceforth referred to as Augbin, was proposed by Wason and Seaman.⁷ We briefly describe this method here, but more details are found in Wason and Seaman.⁷

The Augbin method makes assumptions that the log tumour size ratios follow a multivariate normal distribution, and the probability of new-lesion progression depends only on the observed tumour size at the previous visit. The log tumour size ratios are modelled by

$$(Y_{i1}, Y_{i2})' | z_{i0} \sim N((\mu_{i1}, \mu_{i2})', \Sigma),$$

where $\mu_{i1} = \beta_1 + \omega z_{i0}$, $\mu_{i2} = \beta_2 + \omega z_{i0}$. The new-lesion progression is modelled by using logistic regression models

$$\text{Logit}\{\Pr(D_{i1} = 1 | z_{i0})\} = \alpha_1 + \gamma_1 z_{i0},$$

$$\text{Logit}\{\Pr(D_{i2} = 1 | D_{i1} = 0, z_{i0}, z_{i1})\} = \alpha_2 + \gamma_2 z_{i1}.$$

The probability of response for patient i at time 2 is written by

$$\Pr(S_{i2} = 1 | \theta) = \int_{-\infty}^{c_2} \int_{-\infty}^{\infty} \Pr(D_{i1} = 0 | z_{i0}) \Pr(D_{i2} = 0 | D_{i1} = 0, z_{i0}, z_{i1}) f_{Y_1 Y_2}(y_{i1} y_{i2}; \theta) dy_{i1} dy_{i2},$$

where θ is the vector of parameters from the above models and c_2 is the dichotomisation threshold (usually $\log(0.7)$, representing at least a 30% shrinkage in the tumour size from baseline). The mean response probability is estimated by $\overline{\Pr}(S_2 = 1 | \hat{\theta}) = \frac{1}{n} \sum_{i=1}^n \Pr(S_{i2} = 1 | \hat{\theta})$, where $\hat{\theta}$ is the maximum likelihood estimator of θ . A program is available in the paper which uses R2Cuba to compute the above integration. An approximately $(1-\alpha)\%$ confidence interval for the probability of response is constructed on the logit scale, that is, $\text{expit}\left\{l(\hat{\theta}) \pm \Phi^{-1}\left(1 - \frac{\alpha}{2}\right) \sqrt{\text{var}(l(\hat{\theta}))}\right\}$, where $l(\theta) = \log \frac{\Pr(S_2=1|\theta)}{1-\Pr(S_2=1|\theta)}$ and $\text{var}(l(\hat{\theta}))$ is obtained by using the delta method.

2.4 | Extended Augbin at a fixed time ($t > 2$)

We use the same assumptions and extend the Augbin to t follow-up times. The log tumour size ratios are modelled by

$$(Y_{i1}, \dots, Y_{it})' | z_0 \sim N((\mu_{i1}, \dots, \mu_{it})', \Sigma). \quad (1)$$

An unstructured covariance matrix is used (although an alternative form may be needed if t is large enough). The new-lesion progression is modelled by

$$\text{Logit}\{\Pr(D_{it} = 1 | D_{i1} = \dots = D_{i(t-1)} = 0, z_{i0}, \dots, z_{i(t-1)})\} = \alpha_t + \gamma_t z_{i(t-1)}. \quad (2)$$

We assume that the new-lesion progression depends only on the previous observed tumour size. The missing tumour size because of new-lesion progression can be, therefore, treated as missing at random (MAR) as justified in Wason and Seaman.⁷ $\Pr(Y_{i(t+1)} \text{ is missing} | z_{i0}, \dots, z_{it}) = \Pr(Y_{i(t+1)} \text{ is missing} | z_{i0}, \dots, z_{it})$. See the “sequential missingness at random” section in Supporting Information for details. We also assume that dropout for other reasons before preplanned time is MAR. The probability of response for patient i at time T can be written by

$$\Pr(S_{iT} = 1 | \theta) = \int_{-\infty}^{c_T} \int_{-\infty}^c \dots \int_{-\infty}^c \prod_{t=1}^T \Pr(D_{it} = 0 | D_{i1} = \dots = D_{i(t-1)} = 0, z_{i0}, \dots, z_{i(t-1)}) \times f_{Y_1, \dots, Y_T}(y_{i1}, \dots, y_{iT}; \theta) dy_{i1} \dots dy_{iT}, \quad (3)$$

where c_T and c are the dichotomisation cut points (usually $\log(0.7)$ and $\log(1.2)$, representing at least a 30% shrinkage in the tumour size from baseline and an increase in tumour size by more than 20%). Note that the response at a fixed time $T' \leq T$ can be obtained by replacing T by T' in the above formulas. The advantage is that the Equation 3 uses the models to estimate probability of response of patients and missing data are MAR, it can be applied to patients who drop out or progress before preplanned time. The probability is interpreted as the probability of patient i being a responder at time T as if they were observed until T . A potential issue of Equation 3 is that the multivariate integration is computationally intensive. The mean response probability is estimated by averaging response probability over n patients given $\hat{\theta}$. An approximately $(1 - \alpha)\%$ confidence interval is constructed as described in Section 2.3.

2.5 | Modified Augbin at a fixed time

The objective for this section is to efficiently estimate the mean response probability using continuous tumour-size information in a computationally efficient way. We assume that {no new-lesion progression occurs from time 1 to time T } and {no tumour-growth progression} are conditionally independent given tumour size $\bar{z}_{t-1} = (z_0, \dots, z_{t-1})$. We note this is a strong assumption and assess the sensitivity to this assumption later on. The probability of response for patient i at a fixed time t can be written by

$$\Pr(\text{response} | \bar{z}_{t-1}) = \Pr(\text{no new-lesion progression until } t | \bar{z}_{t-1}) \times \Pr(\text{no tumour progression at time } t | \bar{z}_{t-1}).$$

Let π_t be the probability of new-lesion progression at time t , $t = 1, \dots, T$. Note that π_t is a conditional probability given no new-lesion progression occurring at previous time points. The log tumour size ratio \mathbf{Y}_i is allowed to depend on baseline tumour size whereas new-lesion progression depends on the previous observed tumour size at the previous visit. We can model \mathbf{Y} by

$$(\mathbf{Y}_{i1}, \dots, \mathbf{Y}_{iT})' | z_0 \sim f_{\mathbf{Y}_i}(\cdot), l(\pi_{it}) = \omega_t + \beta_t z_{i(t-1)}. \quad (4)$$

where $f_{\mathbf{Y}_i}(\cdot)$ is a joint distribution and $l(\cdot)$ is the logit link function. We assume that $\pi_{it}(\bar{z}_{t-1}) = \pi_{it}(z_{t-1})$. The probability of response for patient i at a fixed time T can be written by

$$\Pr(S_{iT} = 1 | \bar{z}_{t-1}, \theta) = \prod_{t=1}^T \{1 - \pi_{it}(\bar{z}_{t-1}, \theta)\} \int_{-\infty}^{c_T} \int_{-\infty}^c \dots \int_{-\infty}^c f_{Y_1, \dots, Y_T}(y_{i1}, \dots, y_{iT} | z_{i0}, \theta) dy_{i1} \dots dy_{iT}, \quad (5)$$

where c_T and c are the dichotomisation threshold and θ is a vector of parameters of the models. We assume that $f_{\mathbf{Y}_i}(\cdot)$ is the probability density function of a multivariate normal distribution. The multivariate integration can then be calculated by a highly efficient technique proposed by Genz and Bretz.⁹ The observed data for patient i is (F_i, z_i) . The $l(\pi_{i,F+1})$ can be estimated by $\omega_{F+1} + \beta_{F+1} z_{i,F}$. Their probability of new-lesion progression at time t , $t \geq F_i + 2$ is estimated by

$$\tilde{\pi}_{it} = \frac{1}{n_t - k} \sum_{j: y_{jt} \in \varphi} \pi_{jt}, \quad (6)$$

where n_t is the number of patients with observed z_{t-1} and k is the number of patients who have log tumour size ratio y_t outside of the region of integration φ of Equation 5. We trim those k patients to avoid underestimating π_{it} .

This is similar to an idea of trimmed mean, which is used in many areas and has advantages under both normal and nonnormal distributions.^{10,11}

The vector θ consists of $(T + 1)$ parameters that make up μ , $T(T + 1)/2$ parameters for Σ , and $2T$ parameters for the logistic models. The mean response probability is estimated by $\overline{\Pr}(S_T = 1|\hat{\theta}) = \frac{1}{n} \sum_{i=1}^n \Pr(S_{iT} = 1|\hat{\theta})$, where $\hat{\theta}$ is the maximum likelihood estimator of θ . A $(1 - \alpha)\%$ confidence interval for $\overline{\Pr}(S_T = 1|\theta)$ can be constructed:

$$\left[\overline{\Pr}(S_T = 1|\hat{\theta}) - \Phi^{-1}(1 - \alpha/2) \sqrt{\text{Var}(\overline{\Pr}(S_T = 1|\hat{\theta}))}, \overline{\Pr}(S_T = 1|\hat{\theta}) + \Phi^{-1}(1 - \alpha/2) \sqrt{\text{Var}(\overline{\Pr}(S_T = 1|\hat{\theta}))} \right],$$

where Φ is the standard normal distribution function. However, we found that the method has better properties if we find a confidence interval for $\text{logit}\{\overline{\Pr}(S_T = 1|\hat{\theta})\}$ and transform back. Let $l(\theta) = \log \frac{\Pr(S_T=1|\theta)}{1-\Pr(S_T=1|\theta)}$, we obtain $\text{Var}(l(\hat{\theta}))$ by using the delta method, which is written by

$$\text{var}(l(\hat{\theta})) \approx (\nabla l(\hat{\theta}))^T \text{var}(\hat{\theta}) \nabla l(\hat{\theta}),$$

where $\nabla l(\hat{\theta})$ is the partial derivatives of $l(\theta)$. An approximately $(1 - \alpha)\%$ confidence interval for the probability of response is

$$\left[\text{expit} \left\{ l(\hat{\theta}) - \Phi^{-1}(1 - \alpha/2) \sqrt{\text{var}(l(\hat{\theta}))} \right\}, \text{expit} \left\{ l(\hat{\theta}) + \Phi^{-1}(1 - \alpha/2) \sqrt{\text{var}(l(\hat{\theta}))} \right\} \right].$$

To summarise, the modified method uses a simplification for the relationship between new-lesion progressions and tumour-growth progressions to use a more efficient procedure for multivariate integration.

2.6 | Proposed method for BOR

We focus on the case where confirmation is not required but show briefly how the methodology can straightforwardly allow for confirmation later. By the definition of BOR, a patient is a responder if they have at least 1 log tumour size ratio smaller than $\log(0.7)$ before progression or maximum follow-up time. We define $\Omega_1 = (\log(0.7), \log(1.2))$, $\Omega_2 = (-\infty, \log(0.7))$, and $\Omega_3 = (-\infty, \infty)$ as the possible regions of integration corresponding to being classified as stable disease, responder, and irrelevant variables. Let h be the time at which the patient is first classified as CR/PR. Hence, each component of \bar{Y}_T will fall into 1 of the 3 regions as

$$(Y_1 \dots Y_{h-1} \in \Omega_1, Y_h \in \Omega_2, Y_{h+1} \dots Y_T \in \Omega_3). \quad (7)$$

The probability of response using BOR for patient i will be the sum over all possibilities of when the CR/PR is first observed. Following the concept of the eAugbin, the probability of response can be written by

$$\begin{aligned} \Pr(\text{BOR}_i = 1|\theta) &= \sum_{h=1}^T \int_{\Omega_3^{T-h}} \int_{\Omega_2^1} \int_{\Omega_1^{h-1}} \prod_{t=1}^h \Pr(D_{it} = 0 | D_{i1} = \dots = D_{i(t-1)} = 0, z_{i0}, \dots, z_{i(t-1)}) \\ &\quad \times f_{Y_1, \dots, Y_T}(y_{i1}, \dots, y_{iT}; \theta) dy_{i1} \dots dy_{iT}. \end{aligned} \quad (8)$$

Similarly, following the concept of the mAug, the probability of response can be written by

$$\Pr(\text{BOR}_i = 1|\bar{z}_{t-1}, \theta) = \sum_{h=1}^T \prod_{t=1}^h \{1 - \pi_{it}(\bar{z}_{t-1}, \theta)\} \int_{\Omega_3^{T-h}} \int_{\Omega_2^1} \int_{\Omega_1^{h-1}} f_{Y_1, \dots, Y_T}(y_{i1}, \dots, y_{iT}; z_{i0}, \theta) dy_{i1} \dots dy_{iT}. \quad (9)$$

The mean response probability is then estimated by $\frac{1}{n} \sum_{i=1}^n \Pr(\text{BOR}_i = 1|\hat{\theta})$, where $\hat{\theta}$ is the maximum likelihood estimator of θ . As before, we work on the logit scale, use the delta method to obtain the variance and then transform back to construct the confidence interval for the mean response probability.

When confirmation is required, having 2 continued responses of CR/PR before progression, one can replace (7) with $(Y_1 \dots Y_{h-1} \in \Omega_1, Y_h, Y_{h+1} \in \Omega_2, Y_{h+2} \dots Y_T \in \Omega_3)$ with the sum in (9) going from 1 to $T - 1$.

2.7 | Testing a difference in probability of response between 2 treatments

The above methods can be applied to single-arm trials. For a randomised trial where comparing the difference in response probability is of interest, a minor addition is required.

We assume $2n$ patients are recruited with n patients randomised to each arm. Assumptions for log tumour size ratios and new-lesion progression remain the same as in Section 2.4. We introduce an arm indicator R to the models, with 0 for control and 1 for experimental arms. The log tumour size ratios are modelled by

$$(Y_{i1}, Y_{i2}, \dots, Y_{iT})' | R, z_{i0} \sim N((\mu_{i1}, \mu_{i2}, \dots, \mu_{iT})', \Sigma),$$

where i indexes the i th patient, $\mu_{i1} = \mu_1 + \eta_1 R + \omega z_{i0}$, $\mu_{i2} = \mu_2 + \eta_2 R + \omega z_{i0}$. The new-lesion progression for $T > t$ is modelled by using logistic models

$$\text{Logit}\{\Pr(D_{it} = 1 | D_{i1} = \dots = D_{i(t-1)} = 0, z_{i0}, \dots, z_{i(t-1)})\} = \alpha_t + \beta_t R + \gamma_t z_{i(t-1)}.$$

The probabilities of new-lesion progression for control and experimental arms are $[1 + \exp\{-(\alpha_t + \gamma_t z_{i(t-1)})\}]^{-1}$ and $[1 + \exp\{-(\alpha_t + \beta_t + \gamma_t z_{i(t-1)})\}]^{-1}$, respectively. Let θ be the vector of parameters with length $(T^2 + \frac{9T}{2} + 2)$ from the above models. The mean response probability at a fixed time is estimated by

$$\overline{\Pr}(S = 1 | \hat{\theta}, R) = \frac{1}{2n} \sum_{i=1}^{2n} \Pr(S_i = 1 | \hat{\theta}, R),$$

where $\hat{\theta}$ is the maximum likelihood estimator of θ . We note that patients from both arms are included in the calculation of the probability of response in an arm, as is recommended and justified in Wason and Seaman.⁷ The mean difference in response probability at a fixed time is defined as the difference between mean response probabilities for the 2 arms. It can be written by

$$m_F(\theta) = \overline{\Pr}(S = 1 | R = 1, \theta) - \overline{\Pr}(S = 1 | R = 0, \theta).$$

We obtain the variance of $m_F(\hat{\theta})$ by using the delta method and use the Wald test to test whether $m_F(\theta)$ is 0. Similarly, we define the mean difference in response probability for BOR as

$$m_B(\theta) = \overline{\Pr}(BOR_i = 1 | R = 1, \theta) - \overline{\Pr}(BOR_i = 1 | R = 0, \theta).$$

Both the extended and modified methods can be used as in previous sections. Moreover, we provide a package mAugbin in R including the extended augmented binary method as well as the modified augmented binary method. Changing of integral regions to adapt RECIST or user defined criteria is allowed.

3 | RESULTS

In this section, we evaluate the performance of the proposed methods in terms of precision and power using simulations and a real data example. We use “Bin” to represent the method that just analyses the response outcomes as binary. For single-arm trials, the binary method uses the R-package Hmisc to construct a Wilson interval for binary success ($S = 1$ or $BOR = 1$). For 2-arm studies, the binary method is a logistic regression model that has parameters for treatment group and baseline tumour size, from which the treatment effect can be tested. The terms “Augbin,” “eAugbin,” and “mAug” refer to methods that use continuous information. They are, respectively, Wason and Seaman’s method⁷ at 2 follow-up times, the extended method for more than 2 follow-up times, and the modified method for rapid computation. We use fixed time with varying numbers of follow-up times with c_T and c being $\log(0.7)$ and $-\infty$ and BOR without confirmation as the endpoints.

3.1 | Simulation study setup

Following the aforementioned notation, the observed data available for each patient is (F_i, \mathbf{z}_i) . The observed data are simulated as follows. First of all, baseline tumour size z_{i0} for patient i is generated from a uniform distribution and log tumour size ratios of T follow-up time $\{y_{it} : t = 1, \dots, T\}$ are generated from a multivariate normal distribution. Tumour size z_{it} can then be calculated from $z_{it} = e^{y_{it}} z_{i0}$. Next, new-lesion progression indicators are generated from logistic models with intercept α and tumour size effect γ . A nonzero γ means that probability of new-lesion progression depends on the tumour size at the previous time point. We define time to new-lesion progression as the first time when the new-lesion progression occurs from the logistic models. Finally, tumour size observations of patient i after progression are replaced as missing.

3.1.1 | Single-arm trials assessing response at fixed time

Before generating 5000 replicates, we test the computation time for running one replicate using Augbin/eAugbin. We generated one replicate of 75 patients. Baseline tumour size (Z_0) is generated from a uniform distribution and log

TABLE 1 Mean estimated probability of response and coverage of the modified augmented binary method (mAug) in comparison with using dichotomised continuous method (Bin), augmented binary method (Augbin), and extended augmented binary method (eAugbin) using fixed time with varying numbers of follow-up times

Scenario (α, γ)	Time	True	Mean of estimated probability			Estimated coverage			Reduction in width of 95% CI (%)	
			Bin	Augbin/ eAugbin	mAug	Bin	Augbin/ eAugbin	mAug	Augbin/ eAugbin	mAug
(-1.5, 0)	2	0.334	0.333	0.332	0.338	0.957	0.947	0.947	15.68	14.75
(-2.5, 0.2)	2	0.293	0.293	0.293	0.286	0.948	0.945	0.941	13.45	11.94
(-1.5, 0)	3	0.318	0.316	0.314	0.317	0.953	0.936	0.949	12.53	13.26
(-2.5, 0.2)	3	0.450	0.444	0.443	0.443	0.954	0.943	0.948	14.5	15.28
(-1.5, 0)	4	0.270	0.268	0.263	0.268	0.949	0.926	0.95	12.67	12.54
(-2.5, 0.2)	4	0.429	0.422	0.421	0.421	0.957	0.938	0.943	13.14	14.41

Abbreviation: CI, confidence interval.

tumour size ratios are generated from a multivariate normal distribution for 2 to 6 follow-up times. The (α, γ) are set to $(-1.5, 0)$ and $(-2.5, 0.2)$. The value of $\alpha = -1.5$ corresponds to an 18% chance of developing new lesions between each visit. The computation time for running 1 replicate using Augbin/eAugbin for 2 to 6 follow-up times are 0.04, 0.65, 2.28, 3.41, and 4.47 minutes; while mAug at 6 follow-up times takes 0.09 minutes. We do not consider $T > 4$ because of the length of time need to simulate 5000 replicates for eAugbin. The simulation settings of log tumour size ratios for 2 follow-up times is a similar formulation to Wason and Seaman,⁷ that is,

$$Z_0 \sim U(0, 1), \mathbf{Y}_2 \sim N \left[\log(0.7) \begin{pmatrix} .5 \\ 1 \end{pmatrix}, \begin{pmatrix} .5 & .5 \\ .5 & 1 \end{pmatrix} \right].$$

The settings for $T = 3$ and 4 are $Z_0 \sim U(0, 1)$,

$$\mathbf{Y}_3 \sim N \left[\log(0.7) \begin{pmatrix} .5 \\ .75 \\ 1 \end{pmatrix}, \begin{pmatrix} .5 & .5 & .5 \\ .5 & .75 & .75 \\ .5 & .75 & 1 \end{pmatrix} \right], \mathbf{Y}_4 \sim N \left[\log(0.7) \begin{pmatrix} .25 \\ .5 \\ .75 \\ 1 \end{pmatrix}, \begin{pmatrix} .25 & .25 & .25 & .25 \\ .25 & .5 & .5 & .5 \\ .25 & .5 & .75 & .75 \\ .25 & .5 & .75 & 1 \end{pmatrix} \right].$$

Table 1 shows mean estimated response probability and coverage for Bin, Augbin/eAugbin, and mAug for 2, 3, and 4 follow-up times for 5000 replicates. The columns 10 to 11 show the reduction in 95% confidence interval (CI). They are, respectively, the average of $[1 - (\text{CI width of Augbin})/(\text{CI width of Bin})]$ and $[1 - \text{CI width of mAug}/(\text{CI width of Bin})]$. As seen, in all cases, eAugbin and mAug have narrower CIs compared with Bin. For example, mAug reducing the CI width by 14% means that Bin needs an additional 35% sample size to obtain a similar width. The mAug has a similar coverage to Augbin at $t = 2$. For larger t , it appears the mAug method has a better coverage probability (ie, closer to the nominal value) than eAugbin. The reduction in confidence interval width, compared to the binary method, appears to be similar for the 2 methods. Thus for single-arm trials it appears mAug shows a significant improvement in computational efficiency without notably poorer statistical characteristics compared to eAugbin.

3.1.2 | Randomised trials using response at fixed time

We consider a two-arm trial with a control and experimental arm for 2 follow-up times. Each arm has 75 patients that have been allocated at random. Baseline (Z_0) is generated from a $U(0, 1)$ distribution. The mean log tumour size ratios between each visit are generated from a normal distribution with mean μ and variance $\frac{1}{2}$. We set $\mu = \log(0.7) + \delta\tau + \psi$, where $\delta = 1$ for control and $\delta = -1$ for experimental arms, 2τ is the difference in the mean log tumour size ratio and ψ reflects the effectiveness of the control treatment. This is a similar formulation as.⁷

Figure 1 compares the powers for Bin, eAugbin, and mAug methods for randomised trials. The figure on the right shows the power over treatment effect when $\tau = .35$. As seen, there is a clear power gain when using either mAug or Augbin. mAug performs very closely to Augbin. The empirical Type I error when the difference is 0 for Augbin and mAug are 0.054 and 0.055, respectively.

3.1.3 | Noncomparative trials for BOR

Using the binary composite outcome, patients are classified as responders if they have a CR/PR before time F. The computation time for running 1 replicate using Augbin/eAugbin and BOR for 3 to 6 follow-up times are 0.05, 0.09, 0.3,

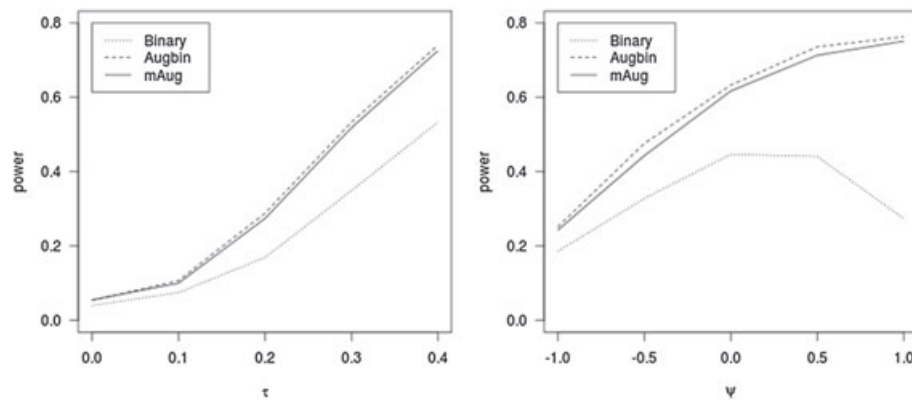


FIGURE 1 Power of the 3 methods for fixed time as the mean log tumour size ratio (τ) varies and as ψ varies at $\tau = .35$

TABLE 2 Mean estimated probability of response and coverage using best observed response without confirmation with Bin, eAugbin, and mAug for maximum number of visits from 4 to 7

(α, γ)	n	Time	True	Mean of estimated probability			Estimated coverage			Reduction in width of 95% CI(%)	
				Bin	eAugbin	mAug	Bin	eAugbin	mAug	eAugbin	mAug
(-1.5, 0)	75	4	0.4	0.4	0.404	0.403	0.959	0.954	0.955	16.5	15.9
(-1.5, 0)	75	5	0.391	0.393	0.398	0.396	0.943	0.951	0.952	16.6	15.9
(-1.5, 0)	75	6	0.386	0.39	0.395	0.394	0.959	0.954	0.955	16.7	16
(-1.5, 0)	150	7	0.382	0.382	—	0.387	0.944	—	0.957	—	16.6
(-2.5, 0.2)	75	4	0.46	0.457	0.462	0.461	0.941	0.957	0.957	16.8	17.3
(-2.5, 0.2)	75	5	0.452	0.448	0.454	0.452	0.954	0.96	0.96	18.2	17.2
(-2.5, 0.2)	75	6	0.446	0.442	0.449	0.447	0.942	0.962	0.961	18.3	17.2
(-2.5, 0.2)	150	7	0.441	0.441	—	0.446	0.95	—	0.96	—	18.3

and 0.56 minutes; while mAug at 6 follow-up times takes 0.22 minutes. Again, we use 5000 replicates of 75 patients. Baseline tumour size is generated from a uniform distribution (0, 1). The log tumour size ratios are generated from multivariate normal distribution for 4, 5, 6, and 7 follow-up times with $\sigma_{tt}^2 = 1$, $t = 4, 5, 6, 7$. Regardless of the number of visits after baseline, we set the mean log tumour size ratios at the end of the treatment to $\log(0.7)$. For example, the case where $T = 4$ refers to having 4 visits after baseline and μ being set to $0.25 \log(0.7)$, $0.5 \log(0.7)$, $0.75 \log(0.7)$, $\log(0.7)$. For computational reasons, eAugbin was included for up to $T = 6$. Table 2 shows the operating characteristics of eAugbin, mAug, and Bin for maximum number of visits varying from 4 to 7. Overall, mAug reduces the average width of the CI by at least 16% compared with Bin. This is equivalent to needing a sample size of around 101 ($1.16^2 \times 75$), to obtain a similar average width using Bin. The reduction in width is slightly higher when there is a tumour size effect on new-lesion progression.

3.1.4 | Comparative trials for BOR

To illustrate results of the mAug method for a two-arm trial, we consider the case where each arm has 75 patients and patients are followed for 4 time points. The mean log tumour size ratios for each time point is $(\log(0.7) + .25\delta\tau)$, $(\log(0.7) + .5\delta\tau)$, $(\log(0.7) + .75\delta\tau)$ and $(\log(0.7) + \delta\tau)$, where $\delta = 1$ for control and $\delta = -1$ for experimental arms respectively. We also consider the method of Karrison et al.⁶ They proposed to assign the highest observed log tumour ratio (referred to as worst outcome henceforth) for deaths and dropouts, and best possible outcome for complete responders, but do not explicitly say how to deal with BOR. We use worst outcome only if the first non-baseline observation was a progression and otherwise use the lowest logtumour before progression. Figure 2 compares the powers for Bin, mAug, and Karrison's methods in comparative trials for 4 time points when BOR is used. The empirical type I error when the difference is 0 for Binary, mAug, and Karrison are 0.041, 0.058, and 0.042, respectively. Although there is a slight inflation in type I error rate for mAug, in general, there is a consistent power advantage when using mAug compared to using Bin and the power of mAug is very similar to Karrison's method.

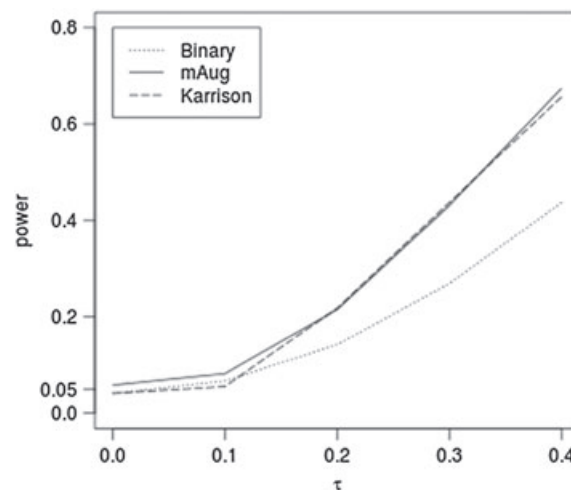


FIGURE 2 Power of the binary, mAug, and Karrison's methods for best observed response for 4 time points as the mean log tumour size ratio (τ) varies

3.2 | Case study: HORIZON II

The HORIZON II (clinicaltrials.gov identifier: NCT00384176) is a 3-arm colon cancer trial sponsored by AstraZeneca. Patients initially were randomly assigned 1:1:1 to placebo, cediranib 20 mg once daily, and cediranib 30 mg once daily. Later, subsequent patients were randomly assigned 1:2 to placebo or cediranib 20 mg.¹² The numbers of patients with baseline record for the 3 arms are 346, 484, and 209. The tumour sizes of patients were measured every 6 weeks up to 24 weeks and then every 12 weeks. Figure S1 in Supporting Information shows a waterfall plot for the individual reduction in tumour size at week 24 from the baseline. There are cases that participants are classified as responders before progression which results in different response estimates between fixed time and BOR.

We used a permutation test to calculate the empirical type I error rate. Data from baseline, 6, 12, 18, and 24 weeks, were used. We simulated 5000 replicates, with the treatment assignment label shuffled randomly in each replicate. For each replicate, we tested the difference in probability of BOR between 2 treatment arms using mAug with 4 follow-up times. The empirical type I error for no difference between placebo and cediranib 20 mg is 0.0558 and that between placebo and cediranib 30 mg is 0.0518. These are within Monte Carlo standard error of a true type I error of 0.05 (MC error ± 0.006).

Figure S2 in Supporting Information shows the mean estimated response probability using the 3 methods and fixed time with between 2 and 5 follow-up times for Placebo, 20 and 30 mg, respectively. The mean estimated response probability decreases as the number of time points increases. Generally, the estimated mean probabilities of response for 3 methods are similar. Figure S3 in Supporting Information shows a residual plot of the fitted multivariate normal model for the 20 mg arm using 3 follow-up times. The residuals look close to normally distributed, though there is a pattern the variance of the residuals may be decreasing as the fitted values increase. In general, it may be beneficial to apply a transformation such as the Box-Cox family.

Table 3 reports the width of the 95% CI for each arms probability of response using fixed time. The width corresponds to the length of the vertical lines shown in Figure S2. The 95% CI widths of eAugbin and mAug are considerably narrower than that of Bin. We compared Placebo and cediranib 20 mg as well as Placebo and cediranib 30 mg using mAug BOR and Bin BOR for 4 to 6 time points. Results show that the mAug method gives a considerably smaller 95% CI than the

TABLE 3 The width of 95% CI for 3 methods using fixed time with between 2 and 5 follow-up times for individual arm

Method	Time Placebo				Cediranib 20 mg				Cediranib 30 mg			
	2	3	4	5	2	3	4	5	2	3	4	5
Bin	0.113	0.114	0.111	0.105	0.111	0.112	0.111	0.111	0.134	0.133	0.13	0.124
eAugbin	0.073	0.074	0.073	0.07	0.072	0.075	0.074	0.064	0.088	0.088	0.088	0.085
mAug	0.086	0.087	0.087	0.08	0.086	0.088	0.086	0.088	0.105	0.105	0.104	0.096

Bin method. The maximum width of the 95% CI for mAug is 0.131 for comparing Placebo with 30 mg, while the width is 0.174 for Bin (See Table S1 in Supporting Information).

4 | DISCUSSION

In this paper, we have considered how the augmented binary method of Wason and Seaman⁷ can be extended to be applicable for a wider range of phase II oncology trials. We have made 3 contributions. The first is to extend the existing method to more than 2 follow-up times. The second is a modified method that considerably reduces the computational time by making a simplifying assumption about the relationship between new lesions and tumour size change. The third is a mechanism for using both of these methods when the endpoint is based around the best observed RECIST observation before progression, which is a common phase II oncology endpoint.

We have shown that all proposed methods carry the same good properties as the augmented binary method. They provide extra precision, ie, they require a smaller sample size for the same precision (compared to the traditional analysis of analysing response as a binary outcome) in single-arm trials and are more powerful in comparative trials. We include Karrison's method, which directly tests the continuous tumour change outcome. As found in Wason and Seaman,⁷ in the comparison of using fixed time, Karrison's method performed better than Augbin when the probabilities of new lesions were different between arms and worse when the mean tumour size changes were different, which we expect to be true for BOR as well. We show results comparing the proposed Augbin on BOR and Karrison's method when the probabilities of new lesions are different between arms. We note in this case the Augbin and Karrison's method give similar power. Karrison's method is simpler to implement, but the Augbin estimates a quantity that is more clinically interpretable.

The difference between the modified (mAug) and extended (eAugbin) methods is that the former uses the estimated probability of new-lesion progression whereas the latter more correctly incorporates variation by averaging all possibilities. Estimation of probabilities using the modified method might be biased if only a few patients remain in a trial at some time point. The mAug has similar properties to eAugbin with respect to precision and power when using BOR.

The extended and modified methods define progression as 20% increase from baseline, whereas RECIST defines progression as 20% increase from the minimum point observed. On the HORIZON II dataset, we examined the number of patients who had their BOR being PR or CR by both of these definitions. The number is the same for both approaches for all number of follow-up times. This indicates that considering progression as being 20% from baseline does not substantially affect the estimation. However, we should point out that the eAugbin would be able to use the RECIST definition of progression by including a suitable indicator variable in the integrand as well as mAug by changing regions of integration of variables. Details about how to use the true definition can be found in Supporting Information.

All proposed augmented binary methods involve modelling the log tumour size ratio and new-lesion progression indicators. The new lesion indicator can include other reasons for progression such as unequivocal progression of nontarget lesions. It may also be possible to include a second logistic regression model for the nontarget lesions separately to increase efficiency further. The log tumour size ratio has been shown to be approximately normally distributed in past data in oncology.¹³ Wason and Seaman⁷ show results from the Augbin can be quite sensitive to deviations from the normal assumption. We find that the residuals of the log tumour ratio in our real data application are close to being normal. However, in general, it may be useful to use a transformation to ensure the normality assumption is as close to true as possible. One could also use other models rather than multivariate normal. An alternative approach is joint latent modelling of longitudinal tumour size data and the new-lesion progression. One can use a random effect model for the repeat tumour size measure and a latent class membership for new-lesion progression. By membership, we mean a participant has probabilities of belonging to latent classes. Each class refers to the time when new-lesion progression occurs. Moreover, tumour-growth progression or new lesions appearing at a time period results in the patient's tumour size measure being missing for all subsequent time periods. Considering this monotone missing pattern in log tumour size, the joint probability of log tumour size can be written as the product of a set of conditional probabilities of current log tumour size ratio given previous data.¹⁴ Complete responses do cause an issue with the proposed methodology as they mean the log tumour ratio is negative infinity—in practice, we would set complete response to an extreme value of log tumour size ratio. Potentially an alternative model such the censored normal distribution¹⁵ could be used. Future work is warranted to investigate whether this more complicated methodology is worth applying.

We have only considered ORR in this work. However, the proposed method can be adjusted for using disease control rate where SD is included. The probability of patients being classified as SD can be calculated by adapting the inner integral in equations from $\log(0.7)$ to $\log(1.2)$. The joint distribution of the probability of being in each ordered category

(CR/PR, SD, and PD) could be estimated by suitably extending the delta-method approach. An increasingly commonly used phase II endpoint is progression-free survival (PFS). Response-based endpoints focus on success, and PFS is interested in progression. Further development of the Augbin from response to progression so that it can be applied to improve analyses of PFS is an area of current work.

5 | SOFTWARE

For estimating the probability of response for fixed time and for BOR, a package mAugbin in R is available at <https://sites.google.com/site/jmswason/supplementary-material> for the methods proposed in this paper. The package includes both the eAugbin and mAug.

ACKNOWLEDGEMENTS

This work was supported by the Medical Research Council (grant number MC_UP_1302/4), Cancer Research UK (grant number C48553/A18113). We thank AstraZeneca for providing HORIZON II data. We thank the editor and two reviewers for their useful comments on improving the paper.

ORCID

Chien-Ju Lin  <http://orcid.org/0000-0002-1440-989X>

REFERENCES

1. Paul SM, Mytelka DS, Dunwiddie CT, et al. How to improve R&D productivity: the pharmaceutical industry's grand challenge. *Nat Rev Drug Discov*. 2010;9:203-214.
2. Johnson JR, Williams G, Pazdur R. End points and United States food and drug administration approval of oncology drugs. *J Clin Oncol*. 2003;21:1404-1411.
3. Pazdur R. Endpoints for assessing drug activity in clinical trials. *Oncologist*. 2008;13:19-21.
4. Eisenhauer EA, Therasse P, Bogaerts J, et al. New response evaluation criteria in solid tumours: revised RECIST guideline (version 1.1). *Eur J Cancer*. 2009;45:228-247.
5. Dhani N, Tu D, Sargent DJ, Seymour L, Moore MJ. Alternate endpoints for screening phase II studies. *Clin Cancer Res*. 2009;15:1873-1882.
6. Karrison TG, Maitland ML, Stadler WM, Ratain MJ. Design of phase II cancer trials using a continuous endpoint of change in tumor size: application to a study of sorafenib and erlotinib in non small-cell lung cancer. *J Natl Cancer Inst*. 2007;99:1455-1461.
7. Wason JMS, Seaman SR. Using continuous data on tumour measurements to improve inference in phase II cancer studies. *Stat Med*. 2013;20:4639-4650.
8. Jaki T, Andre V, Su TL, Whitehead J. Designing exploratory cancer trials using change in tumour size as primary endpoint. *Stat Med*. 2013;32:2544-2554.
9. Genz A, Bretz F. *Computation of Multivariate Normal and t Probabilities*. Heidelberg: Springer-Verlag; 2009.
10. Stigler SM. The asymptotic distribution of the trimmed mean. *Ann Stat*. 1973;1:472-477.
11. Wilcox RR. Trimmed means. In: Everitt BS, Howell D, eds. *Encyclopedia of Statistics in Behavioral Science*. Chichester: John Wiley & Sons; 2005:2066-2067.
12. Hoff PM, Hochhaus A, Pestalozzi BC, et al. Cediranib plus folfox/cafox versus placebo plus folfox/cafox in patients with previously untreated metastatic colorectal cancer: a randomized, double-blind, phase iii study (horizon ii). *J Clin Oncol*. 2012;29:3596-603.
13. Lavin P. An alternative model for the evaluation of antitumour activity. *Cancer Clin Trials*. 1981;4:451-457.
14. Schafer JL. *Analysis of Incomplete Multivariate Data*. London: Chapman & Hall; 1997:218-238.
15. Wason JMS, Mander AP. The choice of test in phase II cancer trials assessing continuous tumour shrinkage when complete responses are expected. *Stat Med*. 2015;24:909-919.

SUPPORTING INFORMATION

Additional Supporting Information may be found online in the supporting information tab for this article.

How to cite this article: Lin C-J, Wason JMS. Improving phase II oncology trials using best observed RECIST response as an endpoint by modelling continuous tumour measurements. *Statistics in Medicine*. 2017;36:4616-4626. <https://doi.org/10.1002/sim.7453>